# NLP methods for crypto price prediction

A. Kipriyanov     G. Kuzmin

ICEF

December 18, 2022
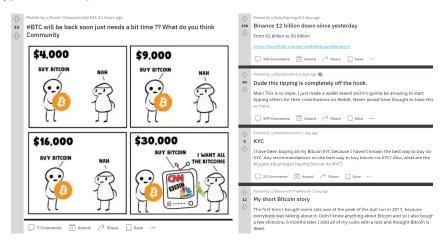
## Idea of the research

- Cryptocurrency – people's money
- Easy access by millions of people
- Social networks (Facebook[1], Twitter, Reddit, etc)
- Pure speculative value creation, no underlying asset
- Retrieve sentiment $\Rightarrow$ forecast the direction of the return
- Compare with more traditional econometric and ML models

---

[1]parent company 'Meta' is recognized as an extremist organization and banned on the territory of the Russian Federation
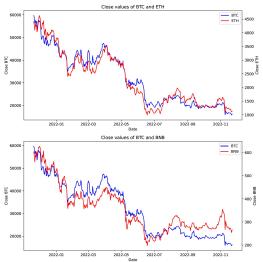
Typical Reddit posts:



We collect Score, Subreddit, Title, Post Text, ID, Total Comments, Date

# Data (2/2)

Yahoo!Finance: Bitcoin, Ethereum, USDT, USDC, Binance Coin and Dogecoin

**Marking Data**

| negative (-1) | neutral (0) | positive (1) |
|---|---|---|
| return $< -2\%$ | return $\in [-2\%; +2\%]$ | return $> +2\%$ |

**TF-IDF**

$$tf = \frac{\# \text{ of the word in the text file}}{\text{total } \# \text{ of all words in the text file}}$$

$$idf = log(\frac{\# \text{ of other text files}}{\# \text{ of text files with this word}})$$

$$tf - idf = tf \times idf$$

- Support Vector Machine (SVM) to perform classification
- Does not count semantic similarity between words

**Word2Vec**

- Attaching semantic numeric vector to each word and obtaining aggregated one for the post
- Finding similar semantic vectors via cosine similarity

**KNN prediction algorithm**

| For post $i = 1$ to $M + V$ ($M$ - train sample, $V$ - test sample): | |
|---|---|
| **Step 1:** for $m = 1$ to $M$: if $i \notin M$, estimate their cosine similarity ($d_{im}$) <br> **Step 2:** Select $KNN$(highest $d$) and take initial sentiments (logreturn) <br> **Step 3:** | **Step 4:** $s_t = \dfrac{\sum\limits_{\{m \in KNN\}} d_{im}^3}{K}$ |
| $Sent_t = \dfrac{\sum\limits_{\{m \in KNN\}} d_{im}^3 \times Initial\ Sent}{\sum\limits_{\{m \in KNN\}} d_{im}^3}$ | **Step 5:** <br> $s_{pred} = \dfrac{\sum\limits_{t=1}^{N} s_t}{N}; \quad Sent_{pred} = \dfrac{\sum\limits_{t=1}^{N} Sent_t \times s_t}{\sum\limits_{t=1}^{N} s_t}$ |

| TF-IDF + SVM | Word2Vec + KNN |
|---|---|
| Train 80%, Test 20% | 2 days to predict next day |
| Expanding window | Rolling window |
| F-beta=0.7027 | F-beta=0.6071 |
| MSE=1.3067 | MSE=1.1429 |

**Limitations and Improvements**

- **TF-IDF**:
  - Add additional factors, such as posts' length, binary for emoji-es and presence of video materials
  - Other alternative classification techniques could be used such as logit, Random Forests
- **Word2Vec**
  - Extremely computationally extensive, that is why, we should check the performance with the training data larger than posts within two days

**Potential applications:**

Constructing trading algorithm, based on predicted marking: (*positive* (1) $\rightarrow$ *Buy*, *neutral* (0) $\rightarrow$ *Hold*, *negative* (-1) $\rightarrow$ *Sell*)

# Competing methodologies (1/2)

- Polynomial regression
    - Lags of BTC, ETH, BNB and DOGE Returns, Lags of BTC Volume and Range, Score and Total Comments
    - Train-validation-test split: 6-2-2
    - CV $\Rightarrow$ power $= 1$
    - MSE $= 0.9185$
- Random Forest
    - Lags of BTC, ETH, BNB and DOGE Returns, Lags of BTC Volume and Range, Score and Total Comments
    - Train-test split: 8-2
    - 10-fold CV $\Rightarrow$ # trees $= 520$
    - **MSE** $= 0.0337$
- GARCH (1,1)
    - Return of BTC
    - Train-test split: 8-2
    - MSE $= 0.0464$

# Competing methodologies (2/2)

- ANN
  - forecasts of Polynomial regression, Random Forest and GARCH (1,1)
  - 2 hidden layers, 4096 neurons each plus a concatenation layer
  - MSE = 0.03507
- ANN+
  - forecasts of all models plus lags of data
  - 2 hidden layers, 4096 neurons each plus a concatenation layer
  - MSE = 1.3806
- RNN
  - Return of BTC
  - Simple RNN layer
  - **MSE** = 0.0248

## Improvements and Further Research

- Increase the data set
- Try more sophisticated ANN (eg. LSTM)
- Parce pictures and determine the sentiment using the body of the post

## Improvements and Further Research

- Increase the data set
- Try more sophisticated ANN (eg. LSTM)
- Parce pictures and determine the sentiment using the body of the post